

1 [0001] SYSTEM AND METHODS OF COOPERATIVELY
2 LOAD-BALANCING CLUSTERED SERVERS

3
4 [0002] Inventors:
5 Tien Le Nguyen
6 Duc Pham
7 Pu Paul Zhang
8 Peter Tsai
9
10

11 [0003] Background of the Invention

12 [0004] Field of the Invention:

13 [0005] The present invention is generally related to systems
14 providing load-balanced network services and, in particular, to techniques for
15 cooperatively distributing load on a cluster of network servers based on
16 interoperation between the cluster of servers and host computers systems that
17 request execution of the network services.

18

19 [0006] Description of the Related Art:

20 [0007] The concept and need for load-balancing arises in a number
21 of different computing circumstances, most often as a requirement for increasing
22 the reliability and scalability of information serving systems. Particularly in the
23 area of networked computing, load-balancing is commonly encountered as a
24 means for efficiently utilizing, in parallel, a large number of information server
25 systems to respond to various processing requests including requests for data from
26 typically remote client computer systems. A logically parallel arrangement of

1 servers adds an intrinsic redundant capability while permitting performance to be
2 scaled linearly, at least theoretically, through the addition of further servers.
3 Efficient distribution of requests and moreover the resulting load then becomes an
4 essential requirement to fully utilizing the paralleled cluster of servers and
5 maximizing performance.

6 [0008] Many different systems have been proposed and variously
7 implemented to perform load-balancing with distinctions typically dependent on
8 the particularities of the load-balancing application. Chung et al. (US Patent
9 6,470,389) describes the use of a server-side central dispatcher that arbitrates the
10 selection of servers to respond to client domain name service (DNS) requests.
11 Clients direct requests to a defined static DNS cluster-server address that
12 corresponds to the central dispatcher. Each request is then redirected by the
13 dispatcher to an available server that can then return the requested information
14 directly to the client. Since each of the DNS requests are atomic and require well-
15 defined server operations, actual load is presumed to be a function of the rate of
16 requests made to each server. The dispatcher therefore implements just a basic
17 hashing function to distribute requests uniformly to the servers participating in the
18 DNS cluster.

19 [0009] The use of a centralized dispatcher for load-balancing control is
20 architecturally problematic. Since all requests flow through the dispatcher, there
21 is an immediate exposure to a single-point failure stopping the entire operation
22 of the server cluster. Further, there is no direct way to scale the performance of
23 the dispatcher. To handle larger request loads or more complex load-balancing
24 algorithms, the dispatcher must be replaced with higher performance hardware
25 at substantially higher cost.

1 [0010] As an alternative, Chung et al. proposes broadcasting all client
2 requests to all servers within the DNS cluster, thereby obviating the need for a
3 centralized dispatcher. The servers implement mutually exclusive hash functions
4 in individualized broadcast request filter routines to select requests for unique local
5 response. This approach has the unfortunate consequence of requiring each
6 server to initially process, to some degree, each DNS request, reducing the
7 effective level of server performance. Further, the selection of requests to service
8 based on a hash of the requesting client address in effect locks individual DNS
9 servers to statically defined groups of clients. The assumption of equal load
10 distribution will therefore be statistically valid, if at all, only over large numbers of
11 requests. The static nature of the policy filter routines also means that all of the
12 routines must be changed every time a server is added or removed from the
13 cluster to ensure that all requests will be selected by a unique server. Given that
14 in a large server cluster, individual server failures are not uncommon and indeed
15 must be planned for, administrative maintenance of such a cluster is likely difficult
16 if not impractical.

17 [0011] Other techniques have been advanced to load-balance networks of
18 servers under various operating conditions. Perhaps the most prevalent load-
19 balancing techniques take the approach of implementing a background or out-of-
20 channel load monitor that accumulates the information necessary to determine
21 when and where to shift resources among the servers dynamically in response to
22 the actual requests being received. For example, Jorden et al. (US Patent
23 6,438,652) describes a cluster of network proxy cache servers where each server
24 further operates as a second level proxy cache for all of the other servers within
25 the cluster. A background load monitor observes the server cluster for repeated

1 second level cache requests for particular content objects. Excessive requests for
2 the same content satisfied from the same second level cache is considered an
3 indication that the responding server is overburdened. Based on a balancing of
4 the direct or first level cache request frequency being served by a server and the
5 second level cache request frequency, the load monitor determines whether to
6 copy the content object to one or more other caches, thereby spreading the
7 second level cache work-load for broadly and repeatedly requested content
8 objects.

9 [0012] Where resources, such as simple content objects, cannot be readily
10 shifted to effect load-balancing, alternate approaches have been developed that
11 characteristically operate by selectively transferring requests, typically represented
12 as tasks or processes, to other servers within a cluster network of servers. Since
13 a centralized load-balancing controller is preferably to be avoided, each server
14 is required to implement a monitoring and communications mechanism to
15 determine which other server can accommodate a request and then actually
16 provide for the corresponding request transfer. The process transfer aspect of the
17 mechanism is often implementation specific in that the mechanism will be highly
18 dependent on the particular nature of the task to transfer and range in complexity
19 from a transfer of a discrete data packet representing the specification of a task
20 to the collection and transport of the entire state of an actively executing process.
21 Conversely, the related conventional load monitoring mechanisms can be
22 generally categorized as source or target oriented. Source oriented servers
23 actively monitor the load status of target servers by actively inquiring of and
24 retrieving the load status of at least some subset of target servers within the
25 cluster. Target oriented load monitoring operates on a publication principle

1 where individual target servers broadcast load status information reflecting, at a
2 minimum, a capacity to receive a task transfer.

3 [0013] In general, the source and target sharing of load status information
4 is performed at intervals to allow other servers within the cluster to obtain on
5 demand or aggregate over time some dynamic representation of the available
6 load capacity of the server cluster. For large server clusters, however, the load
7 determination operations are often restricted to local or server relative network
8 neighborhoods to minimize the number of discrete communications operations
9 imposed on the server cluster as a whole. The trade-off is that more distant server
10 load values must propagate through the network over time and, consequently,
11 result in inaccurate loading reports that lead to uneven distribution of load.

12 [0014] A related problem is described in Allon et al. (US Patent 5,539,883).
13 Server load values, collected into a server cluster load vector, are incrementally
14 requested or advertized by the various servers of the server cluster. Before a
15 server transfers a local copy of the vector, the load values for the server are
16 updated in the vector. Servers receiving the updated vector in turn update the
17 server local copy of the vector with the received load values based on defined
18 rules. Consequently, the redistribution of load values for some given
19 neighborhood may expose an initially lightly loaded server to a protracted high
20 demand for services. The resulting task overload and consequential refusal of
21 service will last at least until a new load vector reflecting the higher server load
22 values circulates among a sufficient number of the servers to properly reflect the
23 load. To alleviate this problem, Allon et al. further describes a tree-structured
24 distribution pattern for load value information as part of the load-balancing
25 mechanism. Based on the tree-structured transfer of load information, low load

1 values, identifying lightly loaded servers, are aged through distribution to preclude
2 lightly loaded servers from being flooded with task transfers.

3 [0015] Whether source or target originated, load-balancing based on the
4 periodic sharing of load information between the servers of the server cluster
5 operates on the fundamental assumption that the load information is reliable as
6 finally delivered. Task transfer rejections are conventionally treated as
7 fundamental failures and, while often recoverable, require extensive exception
8 processing. Consequently, the performance of individual servers may tend to
9 degrade significantly under progressively increasing load, rather than stabilize, as
10 increasing numbers of task transfer recovery and retries operations are required
11 to ultimately achieve a balanced load distribution.

12 [0016] In circumstances where high load conditions are normally incurred,
13 specialized network protocols have been developed to accelerate the exchange
14 and certainty of loading information. Routers and other switch devices are often
15 clustered in various configurations to share network traffic load. A linking network
16 protocol is provided to provide fail-over monitoring in local redundant router
17 configurations and to share load information between both local and remote
18 routers. Current load information, among other shared information, is
19 propagated at high frequency between devices to continuously reflect the
20 individual load status of the clustered devices. As described in Bare (US Patent
21 6,493,318) for example, protocol data packets can be richly detailed with
22 information to define and manage the propagation of the load information and
23 to further detail the load status of individual devices within the cluster. Sequence
24 numbers, hop counts, and various flag-bits are used in support of spanning tree-
25 type information distribution algorithms to control protocol packet propagation

1 and prevent loop-backs. The published load values are defined in terms of
2 internal throughput rate and latency cost, which allows other clustered routers a
3 more refined basis for determining preferred routing paths. While effective, the
4 custom protocol utilized by the devices described in Bare essentially requires that
5 substantial parts of the load-balancing protocol be implemented in specialized,
6 high-speed hardware, such as network processors. The efficient handling of such
7 protocols is therefore limited to specialized, not general purpose computer
8 systems.

9 [0017] Ballard (US Patent 6,078,960) describes a client/server system
10 architecture that, among other features, effects a client-directed load-balanced
11 use of a server network. For circumstances where the various server computer
12 systems available for use by client computer systems may be provided by
13 independent service providers and where use of the different servers may involve
14 different cost structures, Ballard describes a client-based approach for selectively
15 distributing load from the clients to distinct individual servers within the server
16 network. By implementing client-based load-balancing, the client computer
17 systems in Ballard are essentially independent of the service provider server
18 network implementation.

19 [0018] To implement the Ballard load-balancing system, each client
20 computer system is provided with a server identification list from which servers are
21 progressively selected to receive client requests. The list specifies load control
22 parameters, such as the percentage load and maximum frequency of client
23 requests that are to be issued, for each server identified in the list. Server loads
24 are only roughly estimated by the clients based on the connection time necessary
25 for a request to complete or the amount of data transferred in response to a

1 request. Client requests are then issued by the individual clients to the servers
2 selected as necessary to statistically conform to the load-balancing profile defined
3 by the load control parameters. While the server identification list and included
4 load control parameters are static as held by a client, the individual clients may
5 nonetheless retrieve new server identification lists at various intervals from
6 dedicated storage locations on the servers. Updated server identification lists are
7 distributed to the servers as needed under the manual direction of an
8 administrator. Updating of the server identification lists allows an administrator
9 to manually adjust the load-balance profiles as needed due to changing client
10 requirements and to accommodate the addition and removal of servers from the
11 network.

12 [0019] The static nature of the server identification lists makes the client-
13 based load-balancing operation of the Ballard system fundamentally
14 unresponsive to the actual operation of the server network. While specific server
15 loading can be estimated by the various clients, only complete failures to respond
16 to client requests are detectable and then handled only by excluding a non-
17 responsive server from further participation in servicing client requests.
18 Consequently, under dynamically varying loading conditions, the one sided load-
19 balancing performed by the clients can seriously misapprehend the actual loading
20 of the server network and further exclude servers from participation at least until
21 re-enabled through manual administrative intervention. Such blind exclusion of
22 a server from the server network only increases the load on the remaining servers
23 and the likelihood that other servers will, in turn, be excluded from the server
24 network. Constant manual administrative monitoring of the active server network,
25 including the manual updating of server identification lists to re-enable servers

1 and to adjust the collective client balancing of load on the server network, is
2 therefore required. Such administrative maintenance is quite slow, at least relative
3 to how quickly users will perceive occasions of poor performance, and costly to
4 the point of operational impracticality.

5 [0020] From the forgoing discussion, it is evident that an improved system
6 and methods for cooperatively load-balancing a cluster of servers is needed.
7 There is also a further need, not even discussed in the prior art, for cooperatively
8 managing the configuration of a server cluster, not only with respect to the
9 interoperation of the servers as part of the cluster, but further as a server cluster
10 providing a composite service to external client computer systems. Also,
11 unaddressed is any need for security over the information exchanged between the
12 servers within a cluster. As clustered systems become more widely used for
13 security sensitive purposes, diversion of any portion of the cluster operation
14 through the interception of shared information or introduction of a compromised
15 server into the cluster represents an unacceptable risk.

16

17

18 [0021] Summary of the Invention

19 [0022] Thus, a general purpose of the present invention is to provide an
20 efficient system and methods of cooperatively load-balancing a cluster of servers
21 to effectively provide a scalable network service.

22 [0023] This is achieved in the present invention by providing a cluster of
23 servers configured to perform a defined network service. Host computer systems
24 engage in independent transactions with servers of the cluster to distribute
25 requests for the performance of the network service, typically involving a transfer

1 processing of data. The host computer systems are provided with an identification
2 of the servers of the cluster from which the host computer systems dynamically
3 select targeted servers of the cluster with which to conduct respective transactions.
4 The selection of cluster servers is performed autonomously by the host computer
5 systems based on server performance information gathered by host computer
6 systems from individual servers through prior transactions. The cluster server
7 performance information includes load values returned within prior transactions.
8 A returned set of load values reflects the performance status of the corresponding
9 cluster server. Optionally, a concurrently returned weight value reflects a targeted
10 cluster server localized policy evaluation of certain access attribute information
11 provided in conjunction with the service request. A targeted server may explicitly
12 reject a service request based explicitly on the access attributes evaluated locally
13 relative to the operation specified by the network request, load value, weight
14 value, or on a combination thereof. Whether the request is accepted or rejected,
15 the determined load and optional weight values are returned to the request
16 originating host computer to store and use as a basis for selecting a target server
17 for a subsequent transaction.

18 [0024] Thus, an advantage of the present invention is that the necessary
19 operations to effectively load-balance a cluster of server computer systems are
20 cooperatively performed based on autonomous actions implemented between the
21 host computer systems and the targeted servers of the cluster. Load related
22 information is shared in the course of individual service transactions between
23 hosts and cluster servers rather than specifically in advance of individual service
24 transactions. No independent explicit communications connections are required
25 to share loading information among the participating hosts, among the servers

1 of the cluster, or even between the hosts and servers. Consequently, there is no
2 lost performance on the part of the hosts or servers in performing ongoing load-
3 information sharing operations and, moreover, the operational complexity and
4 delay of opening and operating multiple network connections to share loading
5 information is avoided.

6 [0025] Another advantage of the present invention is that the processing
7 overhead incurred to fully utilize the server cluster of the present invention is both
8 minimal and essentially constant relative to service request frequency for both host
9 and server computer systems. Host computer systems perform a substantially
10 constant basis evaluation of available cluster servers in anticipation of issuing a
11 service request and subsequently recording the server response received. Subject
12 to a possible rejection of the request, no further overhead is placed on the host
13 computer systems. Even where a service request rejection occurs, the server
14 selection evaluation is reexecuted with minimal delay or required processing steps.
15 On the server side, each service request is received and evaluated through a
16 policy engine that quickly determines whether the request is to be rejected or, as
17 a matter of policy, given a weight by which to be relatively prioritized in
18 subsequent selection evaluations.

19 [0026] A further advantage of the present invention is that the function of
20 the host computer systems can be distributed in various architectural
21 configurations as needed to best satisfy different implementation requirements.
22 In a conventional client/server configuration, the host function can be
23 implemented directly on clients. Also in a client/server configuration, the host
24 function can be implemented as a filesystem proxy that, by operation of the host,
25 supports virtual mount points that operate to filter access to the data stores of core

1 network file servers. For preferred embodiments of the present invention, the host
2 computer systems are generally the directly protected systems having or providing
3 access to core network data assets.

4 [0027] Still another advantage of the present invention is that the
5 cooperative interoperation of the host systems and the cluster servers enables fully
6 load-balanced redundancy and scalability of operation. A network services cluster
7 can be easily scaled and partitioned as appropriate for maintenance or to address
8 other implementation factors, by modification of the server lists held by the hosts.
9 List modification may be performed through the posting of notices of to the hosts
10 within transactions to mark the presence and withdrawal of servers from the
11 cluster service. Since the server cluster provides a reliable service, the timing of
12 the server list updates are not critical and need not be performed synchronously
13 across the hosts.

14 [0028] Yet another advantage of the present invention is that select
15 elements of the server cluster load-balancing algorithm can be orthogonally
16 executed by the host and server systems. Preferably, discrete servers evaluate
17 instant load and applicable policy information to shape individual transactions.
18 Based on received load and policy weighting information, hosts preferably
19 perform a generally orthogonal traffic shaping evaluation that evolves over
20 multiple transactions and may further consider external factors not directly evident
21 from within a cluster, such as host/server network communications cost and
22 latency. The resulting cooperative load-balancing operation results in an efficient,
23 low-overhead utilization of the host and server performance capacities.

24

1 [0029] Brief Description of the Drawings

2 [0030] Figure 1A is a network diagram illustrating a system environment
3 within which host computer systems directly access network services provided by
4 a server cluster in accordance with a preferred embodiment of the present
5 invention.

6 [0031] Figure 1B is a network diagram illustrating a system environment
7 within which a preferred core network gateway embodiment of the present
8 invention is implemented.

9 [0032] Figure 2 is a detailed block diagram showing the network
10 interconnection between an array of hosts and a cluster of security processor
11 servers constructed in accordance with a preferred embodiment of the present
12 invention.

13 [0033] Figure 3 is a detailed block diagram of a security processor server
14 as constructed in accordance with a preferred embodiment of the present
15 invention.

16 [0034] Figure 4 is a block diagram of a policy enforcement module control
17 process as implemented in a host computer system in accordance with a preferred
18 embodiment of the present invention.

19 [0035] Figure 5 is a simplified block diagram of a security processor server
20 illustrating the load-balancing and policy update functions shared by a server
21 cluster service provider in accordance with a preferred embodiment of the present
22 invention.

23 [0036] Figure 6 is a flow diagram of a transaction process cooperatively
24 performed between a policy enforcement module process and a selected cluster
25 server in accordance with a preferred embodiment of the present invention.

1 [0037] Figure 7A is a flow diagram of a secure cluster server policy update
2 process as performed between the members of a server cluster in accordance with
3 a preferred embodiment of the present invention.

4 [0038] Figure 7B is a block illustration of a secure cluster server policy
5 synchronization message as defined in accordance with a preferred embodiment
6 of the present invention.

7 [0039] Figure 7C is a block illustration of a secure cluster server policy data
8 set transfer message data structure as defined in accordance with a preferred
9 embodiment of the present invention.

10 [0040] Figure 8 is a flow diagram of a process to regenerate a secure
11 cluster server policy data set transfer message in accordance with a preferred
12 embodiment of the present invention.

13 [0041] Figure 9 is a flow diagram illustrating an extended transaction
14 process performed by a host policy enforcement process to account for a version
15 change in the reported secure cluster server policy data set of a cluster server in
16 accordance with a preferred embodiment of the present invention.

17

18

19 [0042] Detailed Description of the Invention

20 [0043] While system architectures have generally followed a client/server
21 paradigm, actual implementations are typically complex and encompass a wide
22 variety of layered network assets. Although architectural generalities are difficult,
23 in all there are fundamentally common requirements of reliability, scalability, and
24 security. As recognized in connection with the present invention, a specific
25 requirement for security commonly exists for at least the core assets, including the

1 server systems and data, of a networked computer system enterprise. The present
2 invention provides for a system and methods of providing a cluster of servers that
3 provide a security service to a variety of hosts established within an enterprise
4 without degrading access to the core assets while maximizing, through efficient
5 load balancing, the utilization of the security server cluster. Those of skill in the
6 art will appreciate that the present invention, while particularly applicable to the
7 implementation of a core network security service, provides fundamentally enables
8 the efficient, load-balanced utilization of a server cluster and, further, enables the
9 efficient and secure administration of the server cluster. As will also be
10 appreciated, in the following detailed description of the preferred embodiments
11 of the present invention, like reference numerals are used to designate like parts
12 as depicted in one ore more of the figures.

13 [0044] A basic and preferred system embodiment 10 of the present
14 invention is shown in Figure 1A. Any number of independent host computer
15 systems 12_{1-N} are redundantly connected through a high-speed switch 16 to a
16 security processor cluster 18. The connections between the host computer systems
17 12_{1-N}, the switch 16 and cluster 18 may use dedicated or shared media and may
18 extend directly or through LAN or WAN connections variously between the host
19 computer systems 12_{1-N}, the switch 16 and cluster 18. In accordance with the
20 preferred embodiments of the present invention, a policy enforcement module
21 (PEM) is implemented on and executed separately by each of the host computer
22 systems 12_{1-N}. Each PEM, as executed, is responsible for selectively routing
23 security related information to the security processor cluster 18 to discretely qualify
24 requested operations by or on behalf of the host computer systems 12_{1-N}. For the
25 preferred embodiments of the present invention, these requests represent a

1 comprehensive combination of authentication, authorization, policy-based
2 permissions and common filesystem related operations. Thus, as appropriate, file
3 data read or written with respect to a data store, generically shown as data store
4 14, is also routed through the security processor cluster 18 by the PEM executed
5 by the corresponding host computer systems 12_{1..N}. Since all of the operations of
6 the PEMs are, in turn, controlled or qualified by the security processor cluster 18,
7 various operations of the host computer systems 12_{1..N} can be securely monitored
8 and qualified.

9 [0045] An alternate enterprise system embodiment 20 of the present
10 invention implementation of the present invention is shown in Figure 1B. An
11 enterprise network system 20 may include a perimeter network 22 interconnecting
12 client computer systems 24_{1..N} through LAN or WAN connections to at least one
13 and, more typically, multiple gateway servers 26_{1..M} that provide access to a core
14 network 28. Core network assets, such as various back-end servers (not shown),
15 SAN and NAS data stores 30, are accessible by the client computer systems 24_{1..N}
16 through the gateway servers 26_{1..M} and core network 28.

17 [0046] In accordance with the preferred embodiments of the present
18 invention, the gateway servers 26_{1..M} may implement both perimeter security with
19 respect to the client computer systems 14_{1..N} and core asset security with respect
20 to the core network 28 and attached network assets 30 within the perimeter
21 established by the gateway servers 26_{1..M}. Furthermore, the gateway servers 26_{1..M}
22 may operate as application servers executing data processing programs on behalf
23 of the client computer systems 24_{1..N}. Nominally, the gateway servers 26_{1..M} are
24 provided in the direct path for the processing of network file requests directed to
25 core network assets. Consequently, the overall performance of the network

1 computer system 10 will directly depend, at least in part, on the operational
2 performance, reliability, and scalability of the gateway servers 26_{1..M}.

3 [0047] In implementing the security service of the gateway servers 26_{1..M},
4 client requests are intercepted by each of the gateway servers 26_{1..M} and redirected
5 through a switch 16 to a security processor cluster 18. The switch 16 may be a
6 high-speed router fabric where the security processor cluster 18 is local to the
7 gateway servers 26_{1..M}. Alternatively, conventional routers may be employed in a
8 redundant configuration to establish backup network connections between the
9 gateway servers 26_{1..M} and security processor cluster 18 through the switch 16.

10 [0048] For both embodiments 10, 20, shown in Figures 1A and 1B, the
11 security processor cluster 18 is preferably implemented as a parallel organized
12 array of server computer systems, each configured to provide a common network
13 service. In the preferred embodiments of the present invention, the provided
14 network service includes a firewall-based filtering of network data packets,
15 including network file data transfer requests, and the selective bidirectional
16 encryption and compression of file data, which is performed in response to
17 qualified network file requests. These network requests may originate directly with
18 the host computer systems 12_{1..N}, client computer systems 14_{1..N}, and gateway
19 servers 16_{1..M} operating as, for example, application servers or in response to
20 requests received by these systems. The detailed implementation and processes
21 carried out by the individual servers of the security processor cluster 18 are
22 described in copending applications Secure Network File Access Control System,
23 Serial Number 10/201,406, Filed July 22, 2002, Logical Access Block Processing
24 Protocol for Transparent Secure File Storage, Serial Number 10/201,409, Filed
25 July 22, 2002, Secure Network File Access Controller Implementing Access

1 Control and Auditing, Serial Number 10/201,358, Filed July 22, 2002, and
2 Secure File System Server Architecture and Methods, Serial Number 10/271,050,
3 Filed October 16, 2002, all of which are assigned to the assignee of the present
4 invention and hereby expressly incorporated by reference.

5 [0049] The interoperation 40 of an array of host computers 12_{1,x} and the
6 security processor cluster 18 is shown in greater detail in Figure 2. For the
7 preferred embodiments of the present invention, the host computers 12_{1,x} are
8 otherwise conventional computer systems variously operating as ordinary host
9 computer systems, whether specifically tasked as client computer systems, network
10 proxies, application servers, and database servers. A PEM component 42_{1,x} is
11 preferably installed and executed on each of the host computers 12_{1,x} to
12 functionally intercept and selectively process network requests directed to any local
13 and core data stores 14, 30. In summary, the PEM components 42_{1,x} selectively
14 forward specific requests in individual transactions to target servers 44_{1,y} within
15 the security processor cluster 18 for policy evaluation and, as appropriate, further
16 servicing to enable completion of the network requests. In forwarding the
17 requests, the PEM components 42_{1,x} preferably operate autonomously.
18 Information regarding the occurrence of a request or the selection of a target
19 server 44_{1,y} within the security processor cluster 18 is not required to be shared
20 between the PEM components 42_{1,x}, particularly on any time-critical basis.
21 Indeed, the PEM components 42_{1,x} have no required notice of the presence or
22 operation of other host computers 12_{1,x} throughout operation of the PEM
23 components 42_{1,x} with respect to the security processor cluster 18.

24 [0050] Preferably, each PEM component 42_{1,x} is initially provided with a list
25 identification of the individual target servers 44_{1,y} within the security processor

1 cluster 18. In response to a network request, a PEM component 42_{1,x} selects a
2 discrete target server 44 for the processing of the request and transmits the
3 request through the IP switch 16 to the selected target server 44. Particularly
4 where the PEM component 42_{1,x} executes in response to a local client process, as
5 occurs in the case of application server and similar embodiments, session and
6 process identifier access attributes associated with the client process are collected
7 and provided with the network request. This operation of the PEM component
8 42_{1,x} is particularly autonomous in that the forwarded network request is
9 preemptively issued to a selected target server 44 with the presumption that the
10 request will be accepted and handled by the designated target server 44.

11 [0051] In accordance with the present invention, a target servers 44_{1,Y} will
12 conditionally accept a network request depending on the current resources
13 available to the target server 44_{1,Y} and a policy evaluation of the access attributes
14 provided with the network request. Lack of adequate processing resources or a
15 policy violation, typically reflecting a policy determined unavailability of a local or
16 core asset against which the request was issued, will result in the refusal of the
17 network request by a target server 44_{1,Y}. Otherwise, the target server 44_{1,Y}
18 accepts the request and performs the required network service.

19 [0052] In response to a network request, irrespective of whether the request
20 is ultimately accepted or rejected, a target server 44_{1,Y} returns load and,
21 optionally, weight information as part of the response to the PEM component
22 42_{1,x} that originated the network request. The load information provides the
23 requesting PEM component 42_{1,x} with a representation of the current data
24 processing load on the target server 44_{1,Y}. The weight information similarly
25 provides the requesting PEM component 42_{1,x} with a current evaluation of the

1 policy determined prioritizing weight for a particular network request, the
2 originating host 12 or gateway server 26 associated with the request, set of access
3 attributes, and the responding target server 44_{1,Y}. Preferably, over the course of
4 numerous network request transactions with the security processor cluster 18, the
5 individual PEM components 42_{1,X} will develop preference profiles for use in
6 identifying the likely best target server 44_{1,Y} to use for handling network requests
7 from specific client computer systems 12_{1,N} and gateway servers 26_{1,M}. While load
8 and weight values reported in individual transactions will age with time and may
9 further vary based on the intricacies of individual policy evaluations, the ongoing
10 active utilization of the host computer systems 12_{1,N} permits the PEM components
11 42_{1,X} to develop and maintain substantially accurate preference profiles that tend
12 to minimize the occurrence of request rejections by individual target servers 44_{1,Y}.
13 The load distribution of network requests is thereby balanced to the degree
14 necessary to maximize the acceptance rate of network request transactions.

15 [0053] As with the operation of the PEM components 42_{1,X}, the operation
16 of the target servers 44_{1,Y} are essentially autonomous with respect to the receipt
17 and processing of individual network requests. In accordance with the preferred
18 embodiments of the present invention, load information is not required to be
19 shared between the target servers 44_{1,Y} within the cluster 18, particularly in the
20 critical time path of responding to network requests. Preferably, the target servers
21 44_{1,Y} uniformly operate to receive any network requests presented and, in
22 acknowledgment of the presented request, identify whether the request is
23 accepted, provide load and optional weight information, and specify at least
24 implicitly the reason for rejecting the request.

1 [0054] While not particularly provided to share load information, a
2 communications link between the individual target servers 44_{1,Y} within the security
3 processor cluster 18 is preferably provided. In the preferred embodiments of the
4 present invention, a cluster local area network 46 is established in the preferred
5 embodiments to allow communication of select cluster management information,
6 specifically presence, configuration, and policy information, to be securely shared
7 among the target servers 44_{1,Y}. The cluster local area network 46
8 communications are protected by using secure sockets layer (SSL) connections and
9 further by use of secure proprietary protocols for the transmission of the
10 management information. Thus, while a separate, physically secure cluster local
11 area network 46 is preferred, the cluster management information may be routed
12 over shared physical networks as necessary to interconnect the target servers 44_{1,Y}
13 of the security processor cluster 18.

14 [0055] Preferably, presence information is transmitted by a broadcast
15 protocol periodically identifying, using encrypted identifiers, the participating
16 target servers 44_{1,Y} of the security processor cluster 18. The security information
17 is preferably transmitted using a lightweight protocol that operates to ensure the
18 integrity of the security processor cluster 18 by precluding rogue or Trojan devices
19 from joining the cluster 18 or compromising the secure configuration of the target
20 servers 44_{1,Y}. A set of configuration policy information is communicated using an
21 additional lightweight protocol that supports controlled propagation of
22 configuration information, including a synchronous update of the policy rules
23 utilized by the individual target servers 44_{1,Y} within the security processor cluster
24 18. Given that the presence information is transmitted at a low frequency relative
25 to the nominal rate of network request processing, and the security and

1 configuration policy information protocols execute only on the administrative
2 reconfiguration of the security processor cluster 18, such as through the addition
3 of target servers 44_{1,Y} and entry of administrative updates to the policy rule sets,
4 the processing overhead imposed on the individual target servers 44_{1,Y} to support
5 intra-cluster communications is negligible and independent of the cluster loading.

6 [0056] A block diagram and flow representation of the software
7 architecture 50 utilized in a preferred embodiment of the present invention is
8 shown in Figure 3. Generally inbound network request transactions are processed
9 through a hardware-based network interface controller that supports routeable
10 communications sessions through the switch 16. These inbound transactions are
11 processed through a first network interface 52, a protocol processor 54, and a
12 second network interface 54, resulting in outbound transactions redirected
13 through the host computers 12_{1,x} to local and core data processing and storage
14 assets 14, 30. The same, separate, or multiple redundant hardware network
15 interface controllers can be implemented in each target server 44_{1,Y} and
16 correspondingly used to carry the inbound and outbound transactions through the
17 switch 16.

18 [0057] Network request data packets variously received by a target server
19 44 from PEM components 42_{1,x}, each operating to initiate corresponding network
20 transactions against local and core network assets 14, 30, are processed through
21 the protocol processor 54 to initially extract selected network and application data
22 packet control information. Preferably, this control information is wrapped in a
23 conventional TCP data packet by the originating PEM component 42_{1,x} for
24 conventional routed transfer to the target server 44_{1,Y}. Alternately, the control
25 information can be encoded as a proprietary RPC data packet. The extracted

1 network control information includes the TCP, IP, and similar networking protocol
2 layer information, while the extracted application information includes access
3 attributes generated or determined by operation of the originating PEM
4 component 42_{1,x} with respect to the particular client processes and context within
5 which the network request is generated. In the preferred embodiments of the
6 present invention, the application information is a collection of access attributes
7 that directly or indirectly identifies the originating host computer, user and
8 domain, application signature or security credentials, and client session and
9 process identifiers, as available, for the host computer 12_{1,N} that originates the
10 network request. The application information preferably further identifies, as
11 available, the status or level of authentication performed to verify the user.
12 Preferably, a PEM component 42_{1,x} automatically collects the application
13 information into a defined data structure that is then encapsulated as a TCP
14 network data packet for transmission to a target server 44_{1,Y}.

15 [0058] Preferably, the network information exposed by operation of the
16 protocol processor 54 is provided to a transaction control processor 58 and both
17 the network and application control information is provided to a policy parser 60.
18 The transaction control processor 58 operates as a state machine that controls the
19 processing of network data packets through the protocol processor 54 and further
20 coordinates the operation of the policy parser in receiving and evaluating the
21 network and application information. The transaction control processor 58 state
22 machine operation controls the detailed examination of individual network data
23 packets to locate the network and application control information and, in
24 accordance with the preferred embodiments of the present invention, selectively
25 control the encryption and compression processing of an enclosed data payload.

1 Network transaction state is also maintained through operation of the transaction
2 control processor 58 state machine. Specifically, the sequences of the network
3 data packets exchanged to implement network file data read and write
4 operations, and other similar transactional operations, are tracked as necessary
5 to maintain the integrity of the transactions while being processed through the
6 protocol processor 54.

7 [0059] In evaluating a network data packet identified by the transaction
8 control processor 58 as an initial network request, the policy parser 60 examines
9 selected elements of the available network and application control information.
10 The policy parser 60 is preferably implemented as a rule-based evaluation engine
11 operating against a configuration policy/key data set stored in a policy/key store
12 62. The rules evaluation preferably implements decision tree logic to determine
13 the level of host computer 12_{1-N} authentication required to enable processing the
14 network file request represented by the network file data packet received, whether
15 that level of authentication has been met, whether the user of a request initiating
16 host computer 12_{1-N} is authorized to access the requested core network assets,
17 and further whether the process and access attributes provided with the network
18 request are adequate to enable access to the specific local or core network
19 resource 14, 30 identified in the network request.

20 [0060] In a preferred embodiment of the present invention, the decision
21 tree logic evaluated in response to a network request to access file data considers
22 user authentication status, user access authorization, and access permissions.
23 Authentication of the user is considered relative to a minimum required
24 authentication level defined in the configuration policy/key data set against a
25 combination of the identified network request core network asset, mount point,

1 target directory and file specification. Authorization of the user against the
2 configuration policy/key data set is considered relative to a combination of the
3 particular network file request, user name and domain, client IP, and client session
4 and client process identifier access attributes. Finally, access permissions are
5 determined by evaluating the user name and domains, mount point, target
6 directory and file specification access attributes with correspondingly specified
7 read/modify/write permission data and other available file related function and
8 access permission constraints as specified in the configuration policy/key data set.

9 [0061] Where PEM components 42_{1,x} function as filesystem proxies, useful
10 to map and redirect filesystem requests for virtually specified data stores to
11 particular local and core network file system data stores 14, 30, data is also
12 stored in the policy/key store 62 to define the set identity of virtual file system
13 mount points accessible to host computer systems 12_{1,N} and the mapping of
14 virtual mount points to real mount points. The policy data can also variously
15 define permitted host computer source IP ranges, whether application
16 authentication is to be enforced as a prerequisite for client access, a limited,
17 permitted set of authenticated digital signatures of authorized applications,
18 whether user session authentication extends to spawned processes or processes
19 with different user name and domain specifications, and other attribute data that
20 can be used to match or otherwise discriminate, in operation of the policy parser
21 60, against application information that can be marshaled on demand by the
22 PEM components 42_{1,x} and network information.

23 [0062] In the preferred embodiments of the present invention, encryption
24 keys are also stored in the policy/key store 62. Preferably, individual encryption
25 keys, as well as applicable compression specifications, are maintained in a

1 logically hierarchical policy set rule structure parseable as a decision tree. Each
2 policy rule provides an specification of some combination of network and
3 application attributes, including the access attributed defined combination of
4 mount point, target directory and file specification, by which permissions
5 constraints on the further processing of the corresponding request can be
6 discriminated. Based on a pending request, a corresponding encryption key is
7 parsed by operation of the policy parser 60 from the policy rule set as needed by
8 the transaction control processor 58 to support the encryption and decryption
9 operations implemented by the protocol processor subject. For the preferred
10 embodiments of the present invention, policy rules and related key data are stored
11 in a hash table permitting rapid evaluation against the network and application
12 information.

13 [0063] Manual administration of the policy data set data is performed
14 through an administration interface 64, preferably accessed over a private
15 network and through a dedicated administration network interface 66. Updates
16 to the policy data set are preferably exchanged autonomously among the target
17 servers 44_{1-Y} of the security processor cluster 18 through the cluster network 46
18 accessible through a separate cluster network interface 68. A cluster policy
19 protocol controller 70 implements the secure protocols for handling presence
20 broadcast messages, ensuring the security of the cluster 46 communications, and
21 exchanging updates to the configuration policy/key data set data.

22 [0064] On receipt of a network request, the transaction control processor
23 58 determines whether to accept or reject the network request dependent on the
24 evaluation performed by the policy parser 60 and the current processing load
25 values determined for the target server 44. A policy parser 60 based rejection will

1 occur where the request fails authentication, authorization, or permissions policy
2 evaluation. For the initially preferred embodiments of the present invention,
3 rejections are not issued for requests received in excess of the current processing
4 capacity of a target server 44. Received requests are buffered and processed in
5 order of receipt with an acceptable increase in the request response latency. The
6 load value immediately returned in response to a request that is buffered will
7 effectively redirect subsequent network requests from the host computers 12_{1,N} to
8 other target servers 44_{1,Y}. Alternately, any returned load value can be biased
9 upward by a small amount to minimize the receipt of network requests that are
10 actually in excess of the current processing capacity of a target server 44. In an
11 alternate embodiment of the present invention, an actual rejection of a network
12 request may be issued by a target server 44_{1,Y} to expressly preclude exceeding the
13 processing capacity of a target server 44_{1,Y}. A threshold of, for example, 95%
14 load capacity can be set to define when subsequent network requests are to be
15 refused.

16 [0065] To provide the returned load value, a combined load value is
17 preferably computed based on a combination of individual load values
18 determined for the network interface controllers connected to the primary network
19 interfaces 52, 56, main processors, and hardware-based encryption/compression
20 coprocessors employed by a target server 44. This combined load value and,
21 optionally, the individual component load values are returned to the request
22 originating host computer 12_{1,N} in response to the network request. Preferably,
23 at least the combined load value is preferably projected to include handling of the
24 current network request. Depending then on the applicable load policy rules

1 governing the operation of the target server 44_{1,Y}, the response returned signals
2 either an acceptance or rejection of the current network request.

3 [0066] In combination with authorization, authentication and permissions
4 evaluation against the network request, the policy parser 60 optionally determines
5 a policy set weighting value for the current transaction, preferably irrespective of
6 whether the network request is to be rejected. This policy determined weighting
7 value represents a numerically-based representation of the appropriateness for
8 use of a particular target server 44 relative to a particular a network request and
9 associated access attributes. For a preferred embodiment of the present
10 invention, a relative low value in a normalized range of 1 to 100, indicating
11 preferred use, is associated with desired combinations of acceptable network and
12 application information. Higher values are returned to identify generally backup
13 or alternative acceptable use. A preclusive value, defined as any value above a
14 defined threshold such as 90, is returned as an implicit signal to a PEM
15 component 42_{1,X} that corresponding network requests are not to be directed to the
16 specific target server 44 except under exigent circumstances.

17 [0067] In response to a network request, a target server 44 returns the reply
18 network data packet including the optional policy determined weighting value, the
19 set of one or more load values, and an identifier indicating the acceptance or
20 rejection of the network request. In accordance with the preferred embodiments
21 of the present invention, the reply network data packet may further specify
22 whether subsequent data packet transfers within the current transaction need be
23 transferred through the security processor cluster 18. Nominally, the data packets
24 of an entire transaction are routed through a corresponding target server 44 to
25 allow for encryption and compression processing. However, where the underlying

1 transported file data is not encrypted or compressed, or where any such
2 encryption or compression is not to be modified, or where the network request
3 does not involve a file data transfer, the current transaction transfer of data need
4 not route the balance of the transaction data packets through the security
5 processor cluster 18. Thus, once the network request of the current transaction
6 has been evaluated and approved by the policy parser 60 of a target server 44,
7 and an acceptance reply packet returned to the host computer 12_{1-N}, the
8 corresponding PEM component 42_{1-X} can selectively bypass use of the security
9 processor cluster 18 for the completion of the current transaction.

10 [0068] An exemplary representation of a PEM component 42, as executed,
11 is shown 80 in Figure 4. A PEM control layer 82, executed to implement the
12 control function of the PEM component 42, is preferably installed on a host system
13 12 as a kernel component under the operating system virtual file system switch or
14 equivalent operating system control structure. In addition to supporting a
15 conventional virtual file system switch interface to the operating system kernel, the
16 PEM control layer 82 preferably implements some combination of a native or
17 network file system or an interface equivalent to the operating system virtual file
18 system switch interface through which to support internal or operating system
19 provided file systems 84. Externally provided file systems 84 preferably include
20 block-oriented interfaces enabling connection to direct access (DAS) and storage
21 network (SAN) data storage assets and file-oriented interfaces permitting access
22 to network attached storage (NAS) network data storage assets.

23 [0069] The PEM control layer 82 preferably also implements an operating
24 system interface that allows the PEM control layer 82 to obtain the hostname or
25 other unique identifier of the host computer system 12, the source session and

1 process identifiers corresponding to the process originating a network file request
2 as received through the virtual file system switch, and any authentication
3 information associated with the user name and domain for the process originating
4 the network file request. In the preferred embodiments of the present invention,
5 these access attributes and the network file request as received by the PEM control
6 layer 82 are placed in a data structure that is wrapped by a conventional TCP
7 data packet. This effectively proprietary TCP data packet is then transmitted
8 through the IP switch 16 to present the network request to a selected target server
9 44. Alternately, a conventional RPC structure could be used in place of the
10 proprietary data structure.

11 [0070] The selection of the target server 44 is performed by the PEM control
12 layer 82 based on configuration and dynamically collected performance
13 information. A security processor IP address list 86 provides the necessary
14 configuration information to identify each of the target servers 44_{1,Y} within the
15 security processor cluster 18. The IP address list 86 can be provided manually
16 through a static initialization of the PEM component 42 or, preferably, is retrieved
17 as part of an initial configuration data set on an initial execution of the PEM
18 control layer 82 from a designated or default target server 44_{1,Y} of the security
19 processor cluster 18. In the preferred embodiment of the present invention, each
20 PEM component 42_{1,X}, in initial execution, implements an authentication
21 transaction against the security processor cluster 18 through which the integrity
22 of the executing PEM control layer 82 is verified and the initial configuration data,
23 including an IP address list 86, is provided to the PEM component 42_{1,X}.

24 [0071] Dynamic information, such as the server load and weight values, is
25 progressively collected by an executing PEM component 42_{1,X} into a SP

1 loads/weights table 88. The load values are timestamped and indexed relative
2 to the reporting target server 44. The weight values are similarly timestamped
3 and indexed. For an initial preferred embodiment, PEM component 42_{1,x} utilizes
4 a round-robin target server 44_{1,y} selection algorithm, where selection of a next
5 target server 44_{1,y} occurs whenever the loading of a current target server 44_{1,y}
6 reaches 100%. Alternately, the load and weight values may be further inversely
7 indexed by any available combination of access attributes including requesting
8 host identifier, user name, domain, session and process identifiers, application
9 identifiers, network file operation requested, core network asset reference, and
10 any mount point, target directory and file specification. Using a hierarchical
11 nearest match algorithm, this stored dynamic information allows a PEM
12 component 42_{1,x} to rapidly establish an ordered list several target servers 44_{1,y}
13 that are both least loaded and most likely to accept a particular network request.
14 Should the first identified target server 44_{1,y} reject the request, the next listed target
15 server 44_{1,y} is tried.

16 [0072] A network latency table 90 is preferably utilized to store dynamic
17 evaluations of network conditions between the PEM control layer 82 and each of
18 the target servers 44_{1,y}. Minimally, the network latency table 90 is used to identify
19 those target servers 44_{1,y} that no longer respond to network requests or are
20 otherwise deemed inaccessible. Such unavailable target servers 44_{1,y} are
21 automatically excluded from the target servers selection process performed by the
22 PEM control layer 82. The network latency table 90 may also be utilized to store
23 timestamped values representing the response latency times and communications
24 cost of the various target servers 44_{1,y}. These values may be evaluated in

1 conjunction with the weight values as part of the process of determining and
2 ordering of the target servers 44_{1,Y} for receipt of new network requests.

3 [0073] Finally, a preferences table 92 may be implemented to provide a
4 default traffic shaping profile individualized for the PEM component 42_{1,X}. For
5 an alternate embodiment of the present invention, a preferences profile may be
6 assigned to each of the PEM components 42_{1,X} to establish a default allocation or
7 partitioning of the target servers 44_{1,Y} within a security processor cluster 18. By
8 assigning target servers 44_{1,Y} different preference values among the PEM
9 components 42_{1,X} and further evaluating these preference values in conjunction
10 with the weight values, the network traffic between the various host computers
11 12_{1,N} and individual target servers 44_{1,Y} can be used to flexibly define use of
12 particular target servers 44_{1,Y}. As with the IP address list 86, the contents of the
13 preferences table may be provided by manual initialization of the PEM control
14 layer 82 or retrieved as configuration data from the security processor cluster 18.

15 [0074] A preferred hardware server system 100 for the target servers 44_{1,Y}
16 is shown in Figure 5. In the preferred embodiments of the present invention, the
17 software architecture 50, as shown in Figure 3, is substantially executed by one
18 or more main processors 102 with support from one or more peripheral,
19 hardware-based encryption/compression engines 104. One or more primary
20 network interface controllers (NICs) 106 provide a hardware interface to the IP
21 switch 16. Other network interface controllers, such as the controller 108,
22 preferably provide separate, redundant network connections to the secure cluster
23 network 46 and to an administrator console (not shown). A heartbeat timer 110
24 preferably provides a one second interval interrupt to the main processors to

1 support maintenance operations including, in particular, the secure cluster
2 network management protocols.

3 [0075] The software architecture 50 is preferably implemented as a server
4 control program 112 loaded in and executed by the main processors 102 from
5 the main memory of the hardware server system 100. In executing the server
6 control program 112, the main processors 102 preferably perform on-demand
7 acquisition of load values for the primary network interface controller 106, main
8 processors 102, and the encryption/compression engines 104. Depending on the
9 specific hardware implementation of the network interface controller 106 and
10 encryption/compression engines 104, individual load values may be read 114
11 from corresponding hardware registers. Alternately, software-based usage
12 accumulators may be implemented through the execution of the server control
13 program 112 by the main processors 102 to track throughput use of the network
14 interface controller 106 and current percentage capacity processing utilization of
15 the encryption/compression engines 104. In the initially preferred embodiments
16 of the present invention, each of the load values represents the percentage
17 utilization of the corresponding hardware resource. The execution of the server
18 control program 112, also provides for establishment of a configuration
19 policy/key data set 116 table also preferably within the main memory of the
20 hardware server system 100 and accessible to the main processors 102. A
21 second table 118 is similarly maintained to receive an updated configuration
22 policy/key data set through operation of the secure cluster network 46 protocols.

23 [0076] Figure 6 provides a process flow diagram illustrating the load-
24 balancing operation 120A implemented by a PEM component 42_{1,X} as executed
25 on a host computer 12_{1,N} cooperatively 120B with a selected target server 44 of

1 the security processor cluster 18. On receipt 122 of a network request from a
2 client 14, typically presented through the virtual filesystem switch to the PEM
3 component 42_{1,x} as a filesystem request, the network request is evaluated by the
4 PEM component 42_{1,x} to associate available access attributes 124, including the
5 unique host identifier 126, with the network request. The PEM component 42_{1,x}
6 then selects 128 the IP address of a target server 44 from the security processor
7 cluster 18.

8 [0077] The proprietary TCP-based network request data packet is then
9 constructed to include the corresponding network request and access attributes.
10 This network request is then transmitted 130 through the IP switch 16 to the target
11 server 44. A target server response timeout period is set concurrently with the
12 transmission 130 of the network request. On the occurrence of a response
13 timeout 132, the specific target server 44 is marked in the network latency table
14 90 as down or otherwise non-responsive 134. Another target server 44 is then
15 selected 128 to receive the network request. Preferably, the selection process is
16 reexecuted subject to the unavailability of the non-responsive target server 44.
17 Alternately, the ordered succession of target servers identified upon initial receipt
18 of the network request may be transiently preserved to support retries in the
19 operation of the PEM component 42_{1,x}. Preservation of the selection list at least
20 until the corresponding network request is accepted by a target server 44 allows
21 a rejected network request to be immediately retried to the next successive target
22 server without incurring the overhead of reexecuting the target server 44 selection
23 process 128. Depending on the duration of the response timeout 132 period,
24 however, re-use of a selection list may be undesirable since any intervening
25 dynamic updates to the security processor loads and weights table 88 and

1 network latency table 90 will not be considered, potentially leading to a higher
2 rate of rejection on retries. Consequently, reexecution of the target server 44
3 selection process 128 taking into account all data in the security processor loads
4 and weights table 88 and network latency table 90 is generally preferred.

5 [0078] On receipt 120B of the TCP-based network request 136, a target
6 server 44 initially examines the network request to access to the request and
7 access attribute information. The policy parser 60 is invoked 138 to produce a
8 policy determined weight value for the request. The load values for the relevant
9 hardware components of the target server 44 are also collected. A determination
10 is then made of whether to accept or reject 140 the network request. If the access
11 rights under the policy evaluated network and application information precludes
12 the requested operation, the network request is rejected. For embodiments of the
13 present invention that do not automatically accept and buffer in all permitted
14 network requests, the network request is rejected if the current load or weight
15 values exceed the configuration established threshold load and weight limits
16 applicable to the target server 44_{1,y}. In either event, a corresponding request
17 reply data packet is generated 142 and returned.

18 [0079] The network request reply is received 144 by the request originating
19 host computer 12_{1,N} and passed directly to the locally executing PEM component
20 42_{1,x}. The load and any returned weight values are timestamped and saved to
21 the security processor loads and weights table 88. Optionally, the network latency
22 between the target server 44 and host computer 12_{1,N}, determined from the
23 network request response data packet, is stored in the network latency table 90.
24 If the network request is rejected 148 based on insufficient access attributes 150,
25 the transaction is correspondingly completed 152 with respect to the host

1 computer 12_{1,N}. If rejected for other reasons, a next target server 44 is selected
2 128. Otherwise, the transaction confirmed by the network request reply is
3 processed through the PEM component 42_{1,X} and, as appropriate, transferring
4 network data packets to the target server 44 as necessary for data payload
5 encryption and compression processing 154. On completion of the client
6 requested network file operation 152, the network request transaction is complete
7 156.

8 [0080] The preferred secure process 160A/160B for distributing presence
9 information and responsively transferring configuration data sets, including the
10 configuration policy/key data, among the target servers 44_{1,Y} of a security
11 processor cluster 18 is generally shown in Figure 7A. In accordance with the
12 preferred embodiments of the present invention, each target server 44 transmits
13 various cluster messages on the secure cluster network 46. Preferably, a cluster
14 message 170, generally structured as shown in Figure 7B, includes a cluster
15 message header 172 that defines a message type, header version number, target
16 server 44_{1,Y} identifier or simply source IP address, sequence number,
17 authentication type, and a checksum. The cluster message header 172 further
18 includes a status value 174 and a current policy version number 146,
19 representing the assigned version number of the most current configuration and
20 configuration policy/key data set held by the target server 44 transmitting the
21 cluster message 170. The status value 174 is preferably used to define the
22 function of the cluster message. The status types include discovery of the set of
23 target servers 44_{1,Y} within the cluster, the joining, leaving and removal of target
24 servers 44_{1,Y} from the cluster, synchronization of the configuration and
25 configuration policy/key data sets held by the target servers 44_{1,Y}, and, where

1 redundant secure cluster networks 46 are available, the switch to a secondary
2 secure cluster network 46.

3 [0081] The cluster message 170, also includes a PK digest 178 that
4 contains a structured list including a secure hash of the public key, the
5 corresponding network IP, and a status field for each target server 44_{1,y} of the
6 security processor cluster 18, as known by the particular target server 44
7 originating a cluster message 170. Preferably, a secure hash algorithm, such as
8 SHA-1, is used to generate the secure public key hashes. The included status field
9 reflects the known operating state of each target server 44, including
10 synchronization in progress, synchronization done, cluster join, and cluster leave
11 states.

12 [0082] Preferably, the cluster message header 172 also includes a digitally
13 signed copy of the source target server 44 identifier as a basis for assuring the
14 validity of a received cluster message 170. Alternately, a digital signature
15 generated from the cluster message header 172 can be appended to the cluster
16 message 170. In either case, a successful decryption and comparison of the
17 source target server 44 identifier or secure hash of the cluster message header
18 172 enables a receiving target server 44 to verify that the cluster message 170 is
19 from a known source target server 44 and, where digitally signed, has not been
20 tampered with.

21 [0083] For the preferred embodiments of the present invention, the target
22 servers 44_{1,y} of a cluster 18 maintain essentially a common configuration to
23 ensure a consistent operating response to any network request made by any host
24 computer 12_{1,x}. To ensure synchronization the configuration of the target servers
25 44_{1,y}, cluster synchronization messages are periodically broadcast 160A on the

1 secure cluster network 46 by each of the target servers 44_{1,Y}, preferably in
2 response to a hardware interrupt generated by the local heartbeat timer 162.
3 Each cluster synchronization message is sent 164 in a cluster message 170 with
4 a synchronization status 174 value, the current policy version level 176 of the
5 cluster 18, and the securely recognizable set of target servers 44_{1,Y} permitted to
6 participate in the security processor cluster 18, specifically from the frame of
7 reference of the target server 44 originating the cluster synchronization message
8 170.

9 [0084] Each target server 44 concurrently processes 160B broadcast cluster
10 synchronization messages 170 as received 180 from each of the other active
11 target servers 44_{1,Y} on the secure cluster network 46. As each cluster
12 synchronization message 170 is received 180 and validated as originating from
13 a target server 44 known to validly exist in the security processor cluster 18, the
14 receiving target server 44 will search 182 the digests of public keys 176 to
15 determine whether the public key of the receiving target server is contained within
16 the digest list 176. If the secure hash equivalent of the public key of a receiving
17 target server 44 is not found 184, the cluster synchronization message 170 is
18 ignored 186. Where the secure hashed public key of the receiving target server
19 44 is found in a received cluster synchronization message 170, the policy version
20 number 174 is compared to the version number of the local configuration
21 policy/key data set held by the receiving target server 44. If the policy version
22 number 174 is the same or less than that of the local configuration policy/key
23 data set, the cluster synchronization message 170 is again ignored 186.

24 [0085] Where the policy version number 174 identified in a cluster
25 synchronization message 170 is greater than that of the current active

1 configuration policy/key data set, the target server 44 issues a retrieval request
2 190, preferably using an HTTPs protocol, to the target server 44 identified within
3 the corresponding network data packet as the source of the cluster
4 synchronization message 170. The comparatively newer configuration policy/key
5 data set held by the identified source target server 44 is retrieved to update the
6 configuration policy/key data set held by the receiving target server 44. The
7 identified source target server 44 responds 192 by returning a source encrypted
8 policy set 200.

9 [0086] As generally detailed in Figure 7C, a source encrypted policy set 200
10 is preferably a defined data structure containing an index 202, a series of
11 encrypted access keys 204_{1,Z}, where Z is the number of target servers 44_{1,Y} known
12 by the identified source target server 44 to be validly participating in security
13 processor cluster 18, an encrypted configuration policy/key data set 206, and a
14 policy set digital signature 208. Since the distribution of configuration policy/key
15 data sets 206 may occur successively among the target servers 44_{1,Y}, the number
16 of valid participating target servers 44_{1,Y} may vary from the viewpoint of different
17 target servers 44_{1,Y} of the security processor cluster 18 while a new configuration
18 policy/key data set version is being distributed.

19 [0087] The index 202 preferably contains a record entry for each of the
20 known validly participating target servers 44_{1,Y}. Each record entry preferably
21 stores a secure hash of the public key and an administratively assigned identifier
22 of a corresponding target server 44_{1,Y}. By convention, the first listed record entry
23 corresponds to the source target server 44 that generated the encrypted policy set
24 200. The encrypted access keys 204_{1,Z} each contain the same triple-DES key,
25 through encrypted with the respective public keys of the known validly

1 participating target servers 44_{1,Y}. The source of the public keys used in encrypting
2 the triple-DES key is the locally held configuration policy/key data set.
3 Consequently, only those target servers 44_{1,Y} that are validly known to the target
4 server 44 that sources an encrypted policy set 200 will be able to first decrypt a
5 corresponding triple-DES encryption key 204_{1,Z} and then successfully decrypt the
6 included configuration policy/key data set 206.

7 [0088] A new triple-DES key is preferably generated using a random
8 function for each policy version of an encrypted policy set 200 constructed by a
9 particular target servers 44_{1,Y}. Alternately, new encrypted policy sets 200 can be
10 reconstructed, each with a different triple-DES key, in response to each HTTPs
11 request received by a particular target servers 44_{1,Y}. The locally held configuration
12 policy/key data set 206 is triple-DES encrypted using the current generated triple-
13 DES key. Finally, a digital signature 208, generated based on a secure hash of
14 the index 202 and list of encrypted access keys 204_{1,Z}, is appended to complete
15 the encrypted policy set 200 structure. The digital signature 208 thus ensures that
16 the source target server 44 identified by the initial secure hash/identifier pair
17 record is in fact the valid source of the encrypted policy set 200.

18 [0089] Referring again to Figure 7A, on retrieval 190 of a source encrypted
19 policy set 200 and further validation as secure and originating from a target
20 server 44 known to validly exist in the security processor cluster 18, the receiving
21 target server 44 searches the public key digest index 202 for digest value
22 matching the public key of the receiving target server 44. Preferably, the index
23 offset location of the matching digest value is used as a pointer to the data
24 structure row containing the corresponding public key encrypted triple-DES key
25 206 and triple-DES encrypted configuration policy/key data set 204. The private

1 key of the receiving target server 44 is then utilized 210 to recover the triple-DES
2 key 206 that is then used to decrypt the configuration policy/key data set 204. As
3 decrypted, the relatively updated configuration policy/key data set 204 is
4 transferred to and held in the update configuration policy/key data set memory
5 118 of the receiving target server 44. Pending installation of the updated
6 configuration policy/key data set 204, a target server 44 holding a pending
7 updated configuration policy/key data set resumes periodic issuance of cluster
8 synchronization messages 170, though using the updated configuration
9 policy/key data set version number 174.

10 [0090] In accordance with the preferred embodiments of the present
11 invention, updated configuration policy/key data sets 204 are relatively
12 synchronously installed as current configuration policy/key data sets 116 to ensure
13 that the active target servers 44_{1,Y} of the security processor cluster 18 are
14 concurrently utilizing the same version of the configuration policy/key data set.
15 Effectively synchronized installation is preferably obtained by having each target
16 server 44 wait 212 to install an updated configuration policy/key data set 204 by
17 monitoring cluster synchronization messages 170 until all such messages contain
18 the same updated configuration policy/key data set version number 174.
19 Preferably, a threshold number of cluster synchronization messages 170 must be
20 received from each active target server 44, defined as those valid target servers
21 44_{1,Z} that have issued a cluster synchronization message 170 within a defined
22 time period, for a target server 44 to conclude to install an updated configuration
23 policy/key data set. For the preferred embodiments of the present invention, the
24 threshold number of cluster synchronization messages 170 is two. From the
25 perspective of each target server 44, as soon as all known active target servers

1 44_{1,Y} are recognized as having the same version configuration policy/key data set,
2 the updated configuration policy/key data set 118 is installed 214 as the current
3 configuration policy/key data set 116. The process 160B of updating of a local
4 configuration policy/key data set is then complete 216.

5 [0091] Referring to Figure 8, an updated configuration policy/key data set
6 is generated 220 ultimately as a result of administrative changes made to any of
7 the information stored as the local configuration policy/key set data.
8 Administrative changes 222 may be made to modify access rights and similar
9 data principally considered in the policy evaluation of network requests. Changes
10 may also be made as a consequence of administrative reconfiguration 224 of the
11 security processor cluster 18, typically due to the addition or removal of a target
12 server 44. In accordance with the preferred embodiments of the present
13 invention, administrative changes 222 are made by an administrator by access
14 through the administration interface 64 on any of the target servers 44_{1,Y}. The
15 administrative changes 222, such as adding, modifying, and deleting policy rules,
16 changing encryption keys for select policy rule sets, adding and removing public
17 keys for known target servers 44, and modifying the target server 44 IP address
18 lists to be distributed to the client computers 12, when made and confirmed by the
19 administrator, are committed to the local copy of the configuration policy/key data
20 set. On committing the changes 222, the version number of the resulting updated
21 configuration policy/key data set is also automatically incremented 226. For the
22 preferred embodiments, the source encrypted configuration policy/key data set
23 200 is then regenerated 228 and held pending transfer requests from other target
24 servers 44_{1,Y}. The cluster synchronization message 170 is also preferably
25 regenerated to contain the new policy version number 174 and corresponding

1 digest set of public keys 176 for broadcast in nominal response to the local
2 heartbeat timer 162. Consequently, the newly updated configuration policy/key
3 data set will be automatically distributed and relatively synchronously installed on
4 all other active target servers 44_{1,Y} of the security processor cluster 18.

5 [0092] A reconfiguration of the security processor cluster 18 requires a
6 corresponding administrative change to the configuration policy/key data set to
7 add or remove a corresponding public key 232. In accordance with the preferred
8 embodiments of the present invention, the integrity of the security processor cluster
9 18 is preserved as against rogue or Trojan target servers 44_{1,Y} by requiring the
10 addition of a public key to a configuration policy/key data set to be made only by
11 a locally authenticated system administrator or through communications with a
12 locally known valid and active target server 44 of the security processor cluster 18.
13 Specifically, cluster messages 170 from target servers 44 not already identified by
14 a corresponding public key in the installed configuration policy/key data set of a
15 receiving target server 44_{1,Y} are ignored. The public key of a new target server 44
16 must be administratively entered 232 on another known and valid target server
17 44 to be, in effect, securely sponsored by that existing member of the security
18 processor cluster 18 in order for the new target server 44 to be recognized.

19 [0093] Consequently, the present invention effectively precludes a rogue
20 target server from self-identifying a new public key to enable the rogue to join the
21 security processor cluster 18. The administration interface 64 on each target
22 server 44 preferably requires a unique, secure administrative login in order to
23 make administrative changes 222, 232 to a local configuration policy/key data
24 set. An intruder attempting to install a rogue or Trojan target server 44 must have
25 both access to and specific security pass codes for an existing active target server

1 44 of the security processor cluster 18 in order to be possibly successful. Since the
2 administrative interface 64 is preferably not physically accessible from the
3 perimeter network 12, core network 18, or cluster network 46, an external breach
4 of the security over the configuration policy/key data set of the security processor
5 cluster 18 is fundamentally precluded.

6 [0094] In accordance with the preferred embodiments of the present
7 intention, the operation of the PEM components 42_{1,x}, on behalf of the host
8 computer systems 12_{1,x}, is also maintained consistent with the version of the
9 configuration policy/key data set installed on each of the target servers 44_{1,y} of
10 the security processor cluster 18. This consistency is maintained to ensure that the
11 policy evaluation of each host computer 12 network request is handled seamlessly
12 irrespective of the particular target server 44 selected to handle the request. As
13 generally shown in Figure 9, the preferred execution 240A of the PEM components
14 42_{1,x} operates to track the current configuration policy/key data set version
15 number. Generally consistent with the PEM component 42_{1,x} execution 120A,
16 following receipt of a network request 122, the last used policy version number
17 held by the PEM component 42_{1,x} is set 242 with the IP address of the selected
18 target server 44, as determined through the target server selection algorithm 128,
19 in the network request data packet. The last used policy version number is set to
20 zero, as is by default the case on initialization of the PEM component 42_{1,x}, to a
21 value based on initializing configuration data provided by a target server 44 of
22 the security processor cluster 18, or to a value developed by the PEM component
23 42_{1,x} through the cooperative interaction with the target servers 44 of the security
24 processor cluster 18. The network request data packet is then sent 130 to the
25 chosen target server 44.

1 [0095] The target server 44 process execution 240B is similarly consistent
2 with the process execution 120B nominally executed by the target servers 44_{1,y}.
3 Following receipt of the network request data packet 136, an additional check
4 244 is executed to compare the policy version number provided in the network
5 request with that of the currently installed configuration policy/key data set. If the
6 version number presented by the network request is less than the installed version
7 number, a bad version number flag is set 246 to force generation of a rejection
8 response 142 further identifying the version number mismatch as a reason for the
9 rejection. Otherwise, the network request is processed consistent with the
10 procedure 120B. Preferably, the target server process execution 240B also
11 provides the policy version number of the locally held configuration policy/key
12 data set in the request reply data packet irrespective of whether a bad version
13 number rejection response 142 is generated.

14 [0096] On receipt 144 specifically of a version number mismatch rejection
15 response, a PEM component 42_{1,x} preferably updates the network latency table
16 90 to mark 248 the corresponding target server 44 as down due to a version
17 number mismatch. Preferably, the reported policy version number is also stored
18 in the network latency table 90. A retry selection 128 of a next target server 44
19 is then performed unless 250 all target servers 44_{1,y} are then determined
20 unavailable based on the combined information stored by the security processor
21 IP address list 86 and network latency table 90. The PEM component 42_{1,x} then
22 assumes 252 the next higher policy version number as received in a bad version
23 number rejection response 142. Subsequent network requests 122 will also be
24 identified 242 with this new policy version number. The target servers 44_{1,y}
25 previously marked down due to version number mismatches are then marked up

1 254 in the network latency table 90. A new target server 44 selection is then
2 made 128 to again retry the network request utilizing the updated policy version
3 number. Consequently, each of the PEM components 42_{1-x} will consistently track
4 changes made to the configuration policy/key data set in use by the security
5 processor cluster 18 and thereby obtain consistent results independent of the
6 particular target server 44 chosen to service any particular network request.

7 [0097] Thus, a system and methods for cooperatively load-balancing a
8 cluster of servers to effectively provide a reliable, scalable network service has
9 been described. While the present invention has been described particularly with
10 reference to a host-based, policy enforcement module inter-operating with a
11 server cluster, the present invention is equally applicable to other specific
12 architectures by employing a host computer system or host proxy to distribute
13 network requests to the servers of a server cluster through cooperative
14 interoperation between the clients and individual servers. Furthermore, while the
15 server cluster service has been described as a security, encryption, and
16 compression service, the system and methods of the present invention are
17 generally applicable to server clusters providing other network services. Also,
18 while the server cluster has been describes as implementing a single, common
19 service, such is only the preferred mode of the present invention. The server
20 cluster may implement multiple independent services that are all cooperatively
21 load-balanced based on the type of network request initially received by a PEM
22 component.

23 [0098] In view of the above description of the preferred embodiments of the
24 present invention, many modifications and variations of the disclosed
25 embodiments will be readily appreciated by those of skill in the art. It is therefore

- 1 to be understood that, within the scope of the appended claims, the invention may
- 2 be practiced otherwise than as specifically described above.